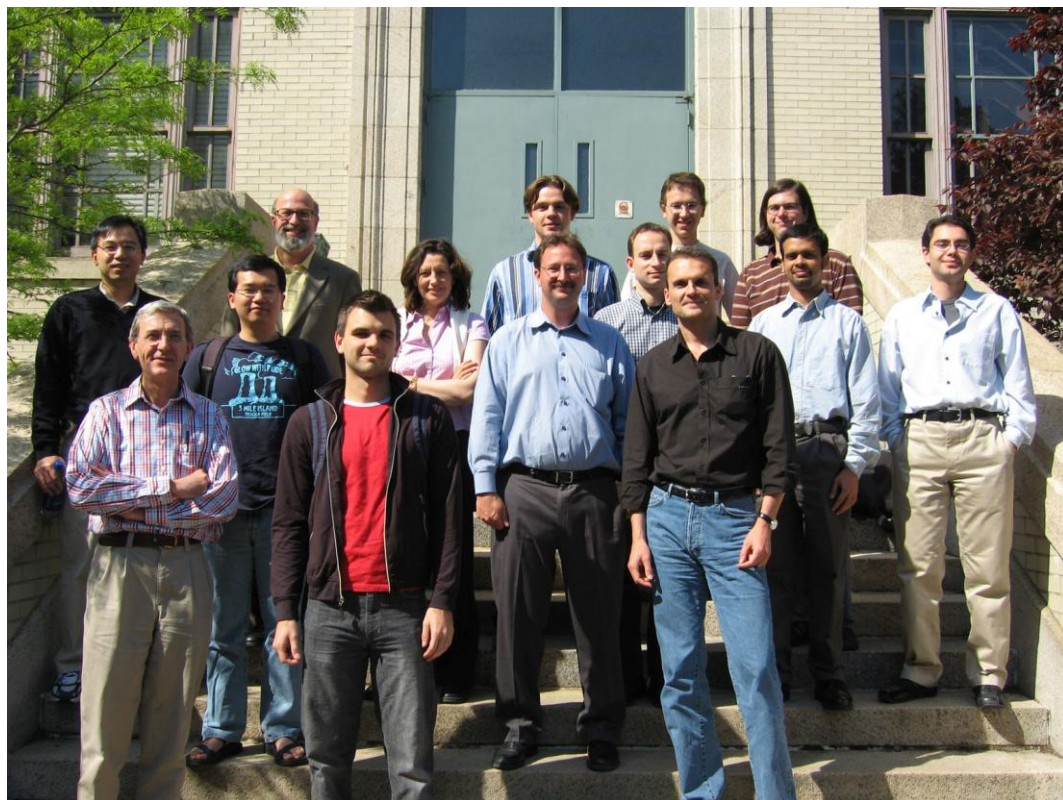**SPIRAL**
www.spiral.net

# Spiral: Program Generation for Linear Transforms and Beyond

**Franz Franchetti**

**Electrical and
Computer Engineering
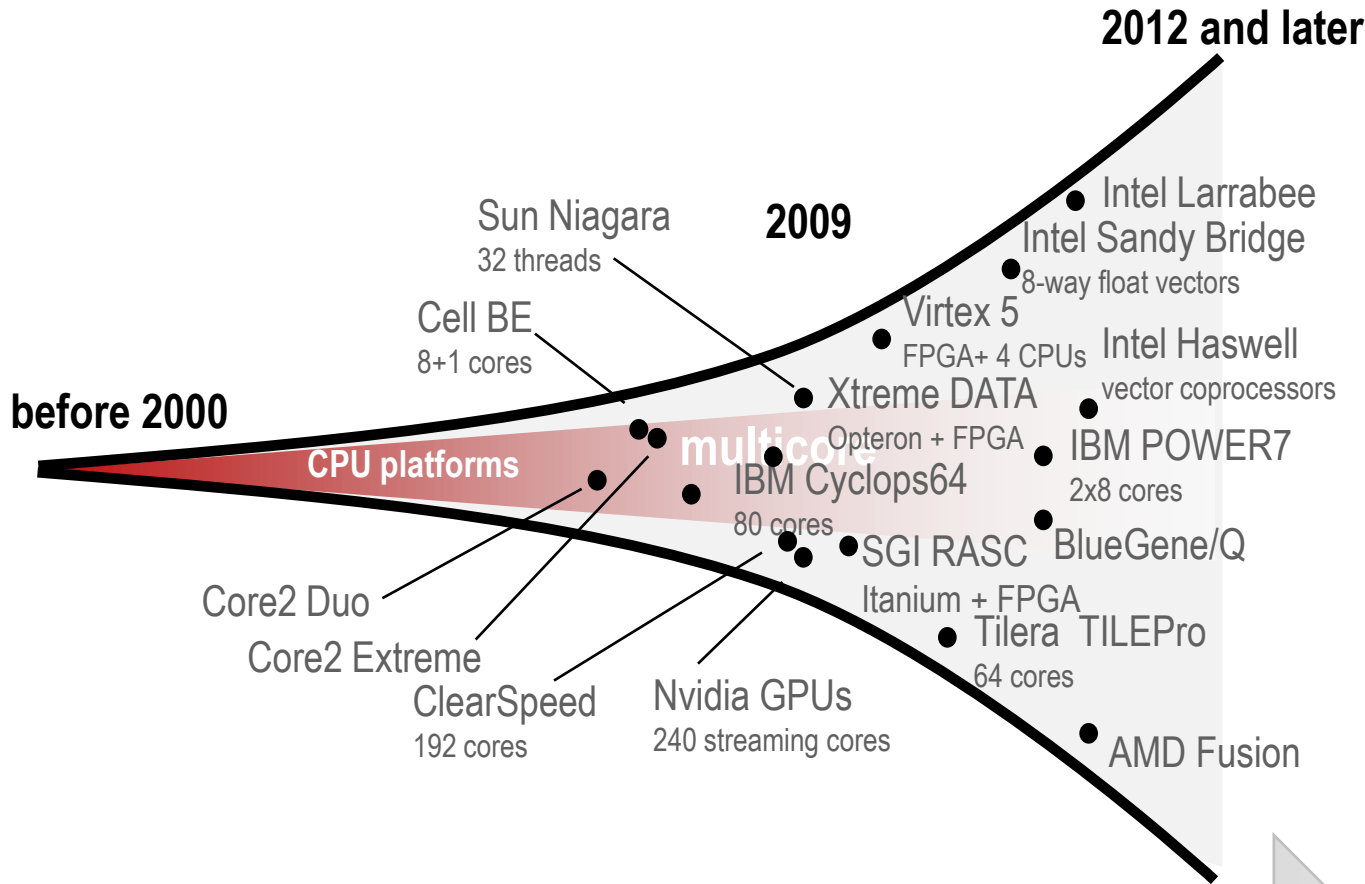Carnegie Mellon University**

**Joint work with**
**Yevgen Voronenko**
**Frédéric de Mesmay**
**Daniel McFarlin**
**Markus Püschel**

**… and the Spiral team (only part shown)**

# The Future is Parallel and Heterogeneous



**2012 and later**

**2009**

Intel Larrabee

Intel Sandy Bridge
8-way float vectors

Sun Niagara
32 threads

Virtex 5
FPGA+ 4 CPUs

Intel Haswell
vector coprocessors

Cell BE
8+1 cores

Xtreme DATA
Opteron + FPGA

IBM POWER7
2x8 cores

**before 2000**

**CPU platforms**

multicore

IBM Cyclops64
80 cores

Core2 Duo

SGI RASC
Itanium + FPGA

BlueGene/Q

Core2 Extreme

Tilera TILEPro
64 cores

ClearSpeed
192 cores

Nvidia GPUs
240 streaming cores

AMD Fusion

?

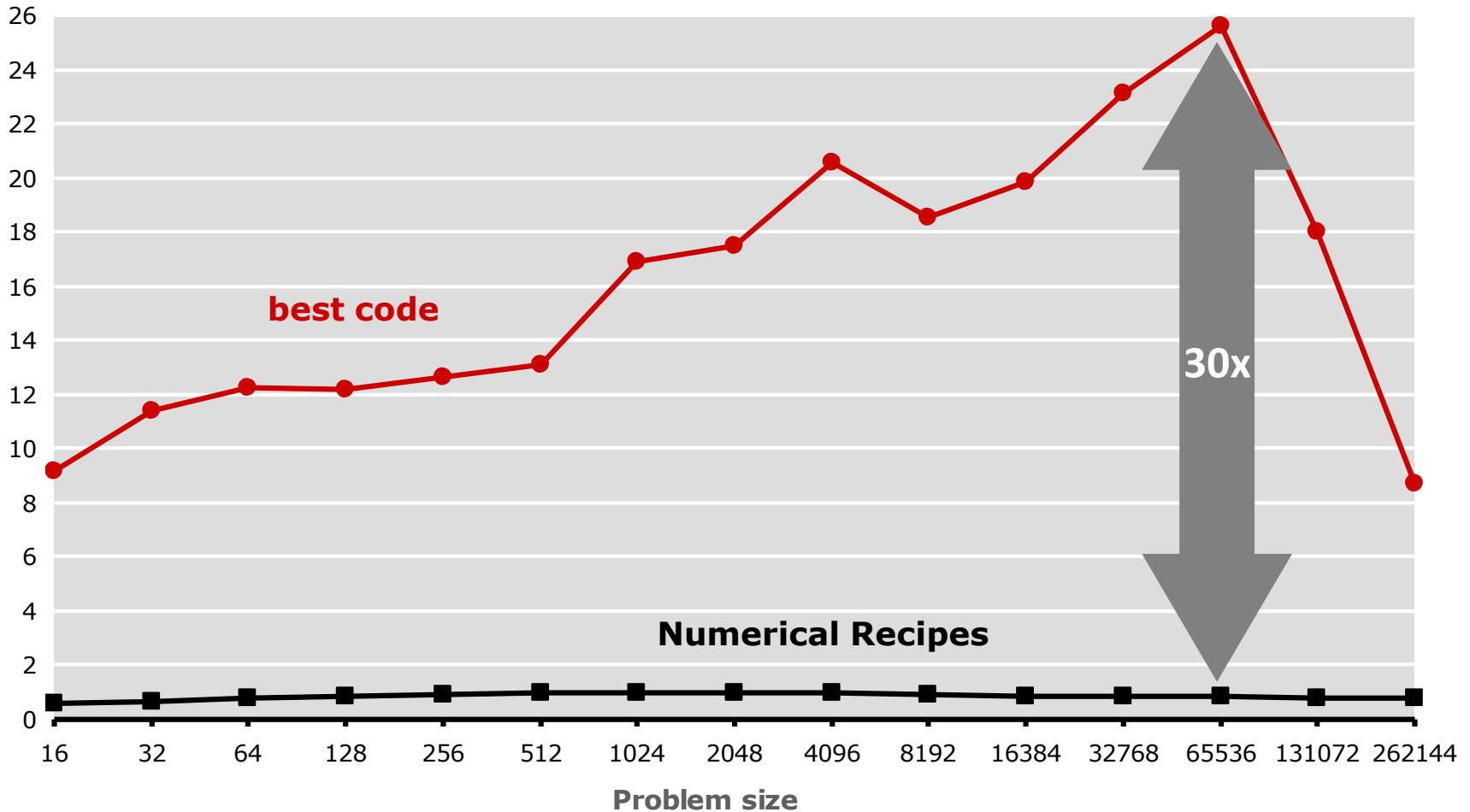*Programmability?*
*Performance portability?*
*Rapid prototyping?*

# The Problem

**Discrete Fourier Transform (single precision): 2 x Core2 Extreme 3 GHz**
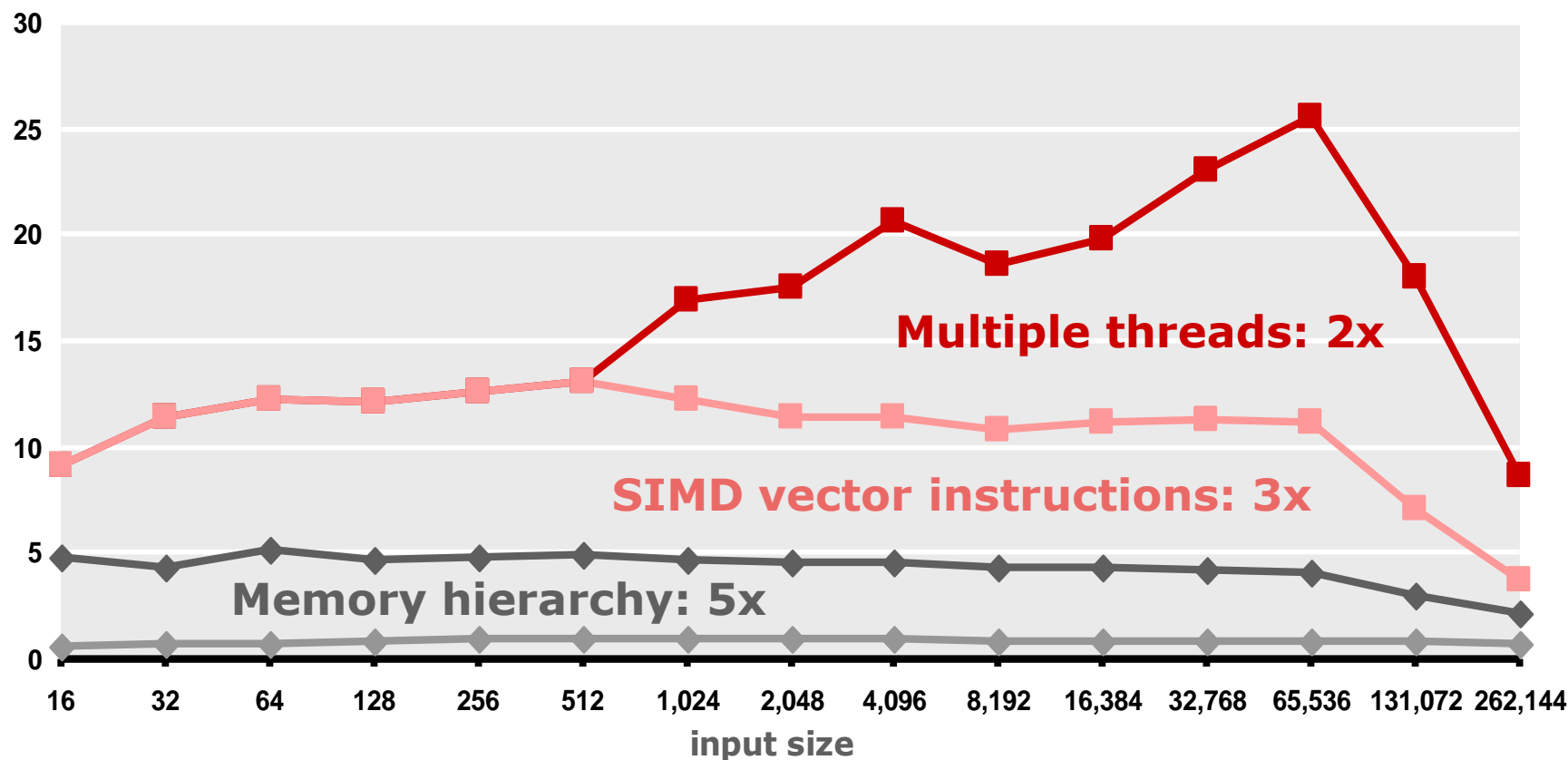Performance [Gflop/s]



*What's going on?*

# DFT Plot: Analysis

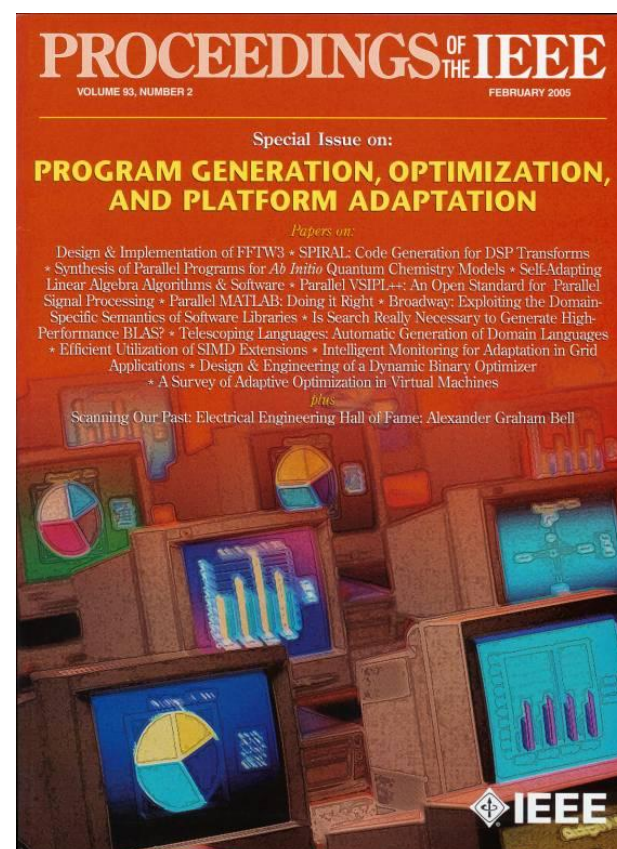**Discrete Fourier Transform (DFT) on 2 x Core 2 Duo 3 GHz**



*High performance library development has become a nightmare*
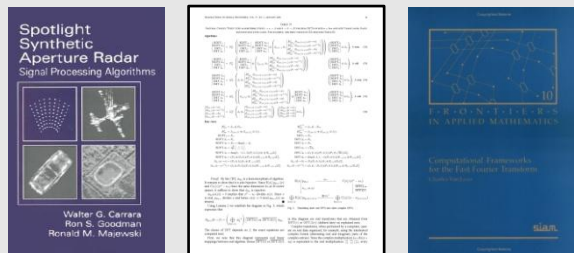
# Automatic Performance Tuning

- **Current vicious circle:** Whenever a new platform comes out, the same functionality needs to be rewritten and reoptimized

- **Automatic Performance Tuning**
  - BLAS: ATLAS, PHiPAC
  - Linear algebra: Sparsity/OSKI, Flame
  - Sorting
  - Fourier transform: FFTW
  - Linear transforms: Spiral
  - …others
  - New compiler techniques
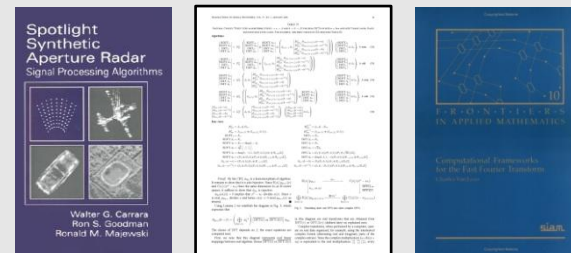
*New challenge: ubiquitous parallelism*



**PROCEEDINGS OF THE IEEE**
VOLUME 93, NUMBER 2 — FEBRUARY 2005

Special Issue on:
**PROGRAM GENERATION, OPTIMIZATION, AND PLATFORM ADAPTATION**

*Papers on:*
Design & Implementation of FFTW3 • SPIRAL: Code Generation for DSP Transforms • Synthesis of Parallel Programs for *Ab Initio* Quantum Chemistry Models • Self-Adapting Linear Algebra Algorithms & Software • Parallel VSIPL++: An Open Standard for Parallel Signal Processing • Parallel MATLAB: Doing it Right • Broadway: Exploiting the Domain-Specific Semantics of Software Libraries • Is Search Really Necessary to Generate High-Performance BLAS? • Telescoping Languages: Automatic Generation of Domain Languages • Efficient Utilization of SIMD Extensions • Intelligent Monitoring for Adaptation in Grid Applications • Design & Engineering of a Dynamic Binary Optimizer • A Survey of Adaptive Optimization in Virtual Machines

*plus*

Scanning Our Past: Electrical Engineering Hall of Fame: Alexander Graham Bell

**IEEE**

**Proceedings of the IEEE special issue, Feb. 2005**

# What is Spiral?

*Traditionally*

*Spiral Approach*



**Spiral**

**High performance library optimized for given platform**

*Comparable performance*

**High performance library optimized for given platform**

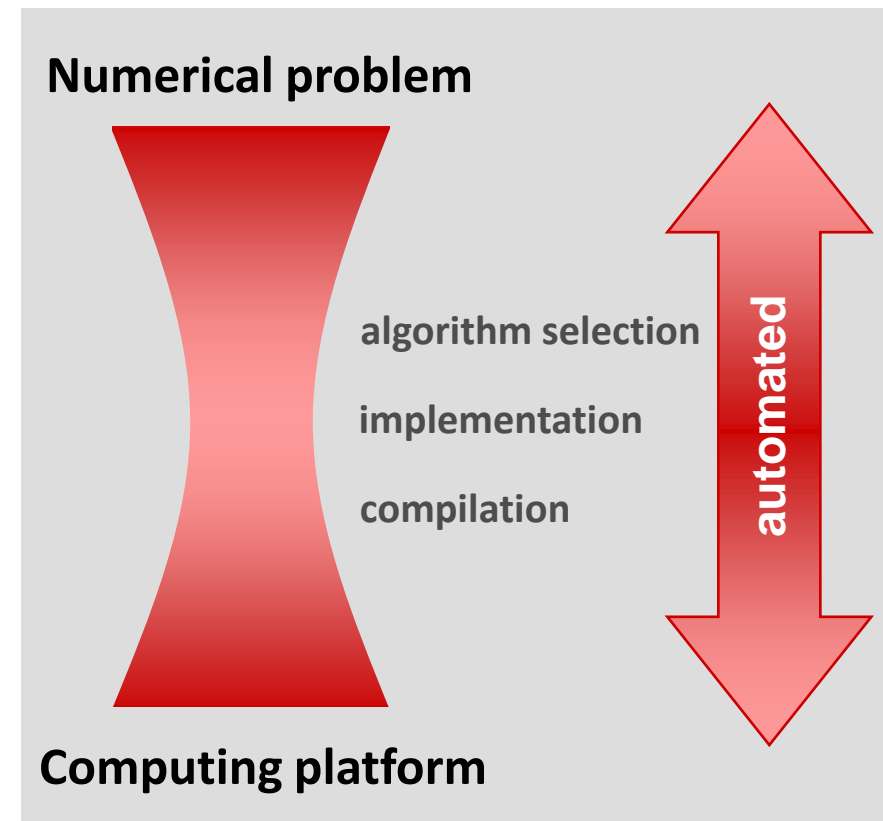# Vision Behind Spiral



*Current*

Numerical problem

C program

algorithm selection

implementation

human effort

automated

compilation

Computing platform

*Future*

Numerical problem

algorithm selection

implementation

compilation

automated

Computing platform

- **C code a singularity: Compiler has no access to high level information**
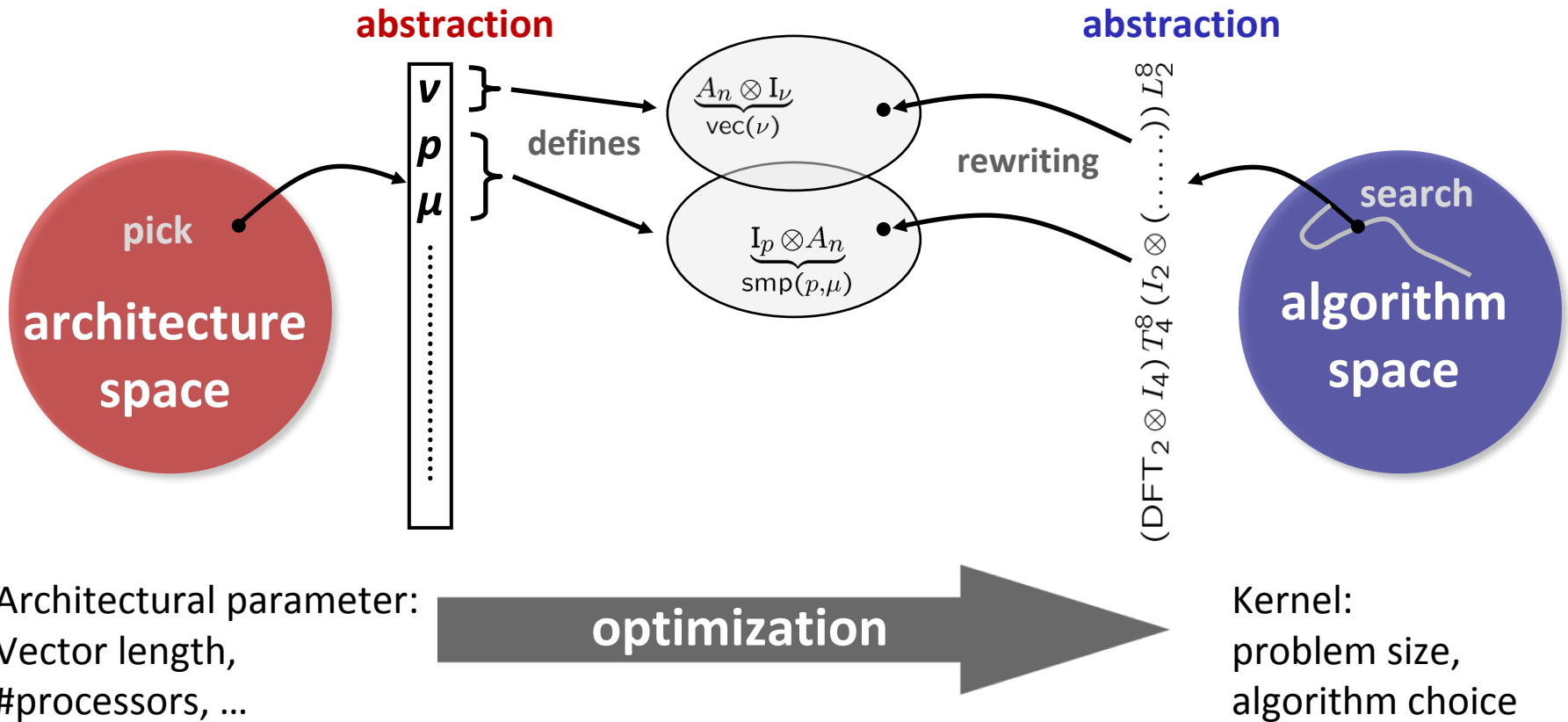
- **Challenge: conquer the high abstraction level for complete automation**

# Main Idea: Program Generation

**Model:** common abstraction
= spaces of matching formulas



Architectural parameter:
Vector length,
#processors, …

**optimization**

Kernel:
problem size,
algorithm choice

# Organization

- **Spiral overview**

- **Parallelization in Spiral**

- **Beyond Transforms**

- **Generating general-size libraries**

- **Results**

- **Concluding remarks**

Markus Püschel, José M. F. Moura, Jeremy Johnson, David Padua, Manuela Veloso, Bryan Singer, Jianxin Xiong, Franz Franchetti, Aca Gacic, Yevgen Voronenko, Kang Chen, Robert W. Johnson, and Nick Rizzolo:
**SPIRAL: Code Generation for DSP Transforms.** Special issue, Proceedings of the IEEE 93(2), 2005

Franz Franchetti, Yevgen Voronenko, Markus Püschel:
**Loop Merging for Signal Transforms.** In Proceedings of Programming Language Design and Implementation (PLDI) 2005.
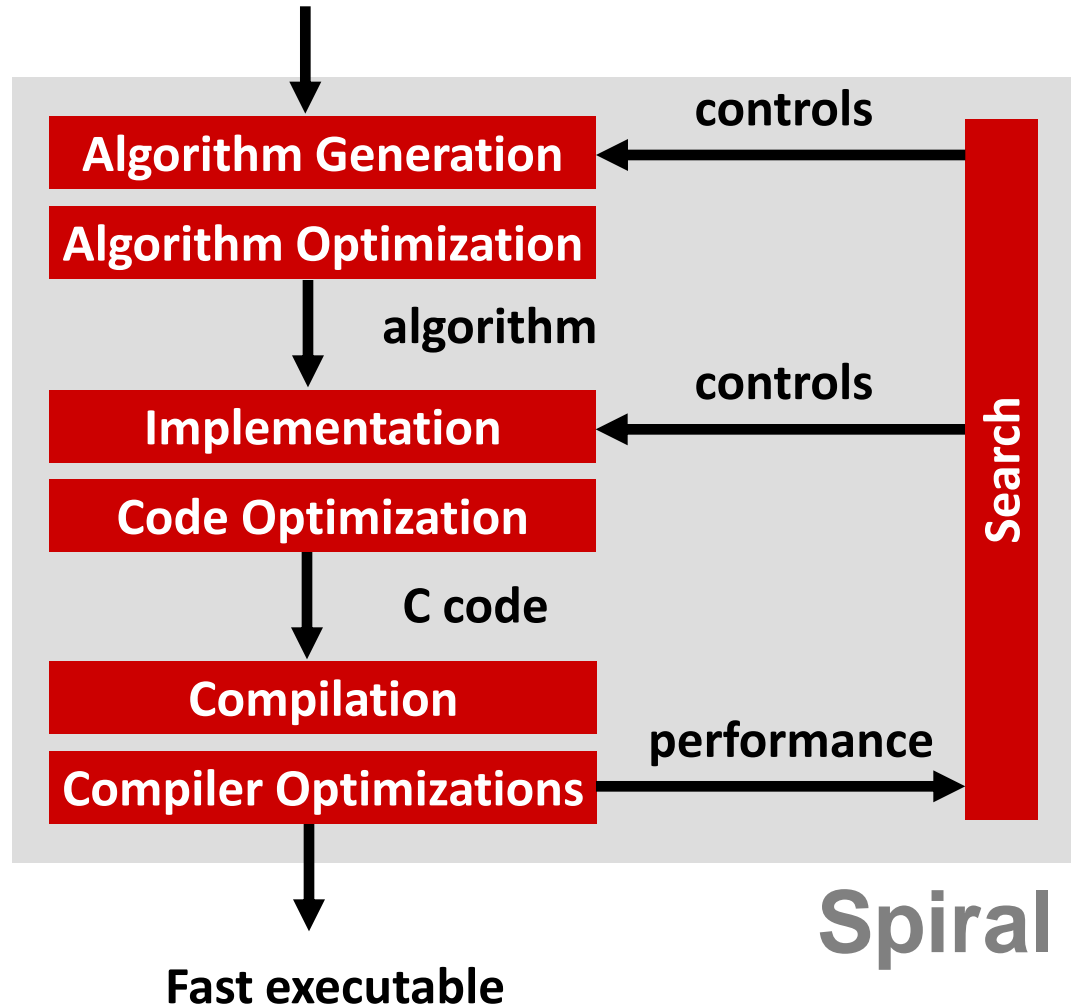
# Spiral

- **Library generator for linear transforms (DFT, DCT, DWT, filters, ….)** *and recently more …*

- **Wide range of platforms supported: scalar, fixed point, vector, parallel, Verilog, GPU**

- **Research Goal: "Teach" computers to write fast libraries**
  - Complete automation of implementation and optimization
  - Conquer the "high" algorithm level for automation

- **When a new platform comes out: Regenerate a retuned library**

- **When a new platform paradigm comes out (e.g., CPU+GPU): Update the tool rather than rewriting the library**

*Intel uses Spiral to generate parts of their MKL and IPP libraries*

# How Spiral Works

**Spiral:**

**Complete automation of the implementation and optimization task**

**Basic idea:**

**Declarative representation of algorithms**

**Rewriting systems to generate and optimize algorithms**

**Problem specification (transform)**

**Algorithm Generation** ← **controls**

**Algorithm Optimization**

**algorithm**

**Implementation** ← **controls**

**Code Optimization**

**C code**

**Compilation**

**Compiler Optimizations** → **performance**

**Search**

**Fast executable**

Spiral

# What is a (Linear) Transform?

■ **Mathematically: Matrix-vector multiplication**

$$x \mapsto y = T \cdot x$$

input vector (signal)

output vector (signal)

<span style="color:red">**transform = matrix**</span>

■ **Example: Discrete Fourier transform (DFT)**

$$\mathbf{DFT}_n \;=\; [e^{-2k\ell\pi i/n}]_{0 \le k,\ell < n}$$

# Transform Algorithms: Example 4-point FFT

**Cooley/Tukey fast Fourier transform (FFT):**

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & j & -1 & -j \\ 1 & -1 & 1 & -1 \\ 1 & -j & -1 & j \end{bmatrix} = \begin{bmatrix} 1 & \cdot & 1 & \cdot \\ \cdot & 1 & \cdot & 1 \\ 1 & \cdot & -1 & \cdot \\ \cdot & 1 & \cdot & -1 \end{bmatrix}\begin{bmatrix} 1 & \cdot & \cdot & \cdot \\ \cdot & 1 & \cdot & \cdot \\ \cdot & \cdot & 1 & \cdot \\ \cdot & \cdot & \cdot & j \end{bmatrix}\begin{bmatrix} 1 & 1 & \cdot & \cdot \\ 1 & -1 & \cdot & \cdot \\ \cdot & \cdot & 1 & 1 \\ \cdot & \cdot & 1 & -1 \end{bmatrix}\begin{bmatrix} 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & 1 & \cdot \\ \cdot & 1 & \cdot & \cdot \\ \cdot & \cdot & \cdot & 1 \end{bmatrix}$$

**Fourier transform**        **Diagonal matrix (twiddles)**

$$\mathrm{DFT}_4 = (\mathrm{DFT}_2 \otimes \mathrm{I}_2)\, \mathsf{T}^4_2 (\mathrm{I}_2 \otimes \mathrm{DFT}_2)\, \mathsf{L}^4_2$$

**Kronecker product**    **Identity**        **Permutation**

- Algorithms reduce arithmetic cost O(n²) → O(nlog(n))
- Product of structured sparse matrices
- Mathematical notation exhibits structure: **SPL (signal processing language)**

# Examples: Transforms

$$\text{DCT-2}_n = \left[\cos(k(2\ell+1)\pi/2n)\right]_{0\le k,\ell<n},$$

$$\text{DCT-3}_n = \text{DCT-2}_n^T \quad \text{(transpose)},$$

$$\text{DCT-4}_n = \left[\cos((2k+1)(2\ell+1)\pi/4n)\right]_{0\le k,\ell<n},$$

$$\text{IMDCT}_n = \left[\cos((2k+1)(2\ell+1+n)\pi/4n)\right]_{0\le k<2n,0\le\ell<n},$$

$$\text{RDFT}_n = [r_{k\ell}]_{0\le k,\ell<n}, \quad r_{k\ell} = \begin{cases} \cos\frac{2\pi k\ell}{n}, & k \le \lfloor\frac{n}{2}\rfloor \\ -\sin\frac{2\pi k\ell}{n}, & k > \lfloor\frac{n}{2}\rfloor \end{cases},$$

$$\text{WHT}_n = \begin{bmatrix} \text{WHT}_{n/2} & \text{WHT}_{n/2} \\ \text{WHT}_{n/2} & -\text{WHT}_{n/2} \end{bmatrix}, \quad \text{WHT}_2 = \text{DFT}_2,$$

$$\text{DHT} = \left[\cos(2k\ell\pi/n) + \sin(2k\ell\pi/n)\right]_{0\le k,\ell<n}.$$

*Spiral currently contains 55 transforms*

# Examples: Breakdown Rules (currently ≈220)

$$\mathrm{DFT}_n \;\rightarrow\; (\mathrm{DFT}_k \otimes \mathrm{I}_m)\, \mathsf{T}^n_m (\mathrm{I}_k \otimes \mathrm{DFT}_m)\, \mathsf{L}^n_k, \quad n = km$$

$$\mathrm{DFT}_n \;\rightarrow\; P_n (\mathrm{DFT}_k \otimes \mathrm{DFT}_m) Q_n, \quad n = km, \; \gcd(k,m) = 1$$

DFT

DCT-3

DCT-4

IMDCT$_{2}$      DCT-4$_{2m}$

- **"Teaches" Spiral about existing algorithm knowledge (~200 journal papers)**

- **Includes many new ones (algebraic theory, Pueschel, Moura, Voronenko)**

$$\mathrm{WHT}_{2^k} \;\rightarrow\; \prod_{i=1}^{t} (\mathrm{I}_{2^{k_1+\cdots+k_{i-1}}} \otimes \mathrm{WHT}_{2^{k_i}} \otimes \mathrm{I}_{2^{k_{i+1}+\cdots+k_t}}), \quad k = k_1 + \cdots + k_t$$

$$\mathrm{DFT}_2 \;\rightarrow\; \mathsf{F}_2$$

$$\mathrm{DCT\text{-}2}_2 \;\rightarrow\; \mathrm{diag}(1, 1/\sqrt{2})\, \mathsf{F}_2$$

$$\mathrm{DCT\text{-}4}_2 \;\rightarrow\; \mathsf{J}_2\, \mathsf{R}_{13\pi/8}$$

**Base case rules**

# SPL to Sequential Code

| SPL construct | code |
|---|---|
| $y = (A_n B_n)x$ | ```t[0:1:n-1] = B(x[0:1:n-1]);```<br>```y[0:1:n-1] = A(t[0:1:n-1];)``` |
| $y = (I_m \otimes A_n)x$ | ```for (i=0;i<m;i++)```<br>```   y[i*n:1:i*n+n-1] =```<br>```      A(x[i*n:1:i*n+n-1])``` |
| $y = (A_m \otimes I_n)x$ | ```for (i=0;i<m;i++)```<br>```   y[i:n:i+m-1] =```<br>```      A(x[i:n:i+m-1]);``` |
| $y = \left( \bigoplus_{i=0}^{m-1} A_n^i \right) x$ | ```for (i=0;i<m;i++)```<br>```   y[i*n:1:i*n+n-1] =```<br>```      A(i, x[i*n:1:i*n+n-1]);``` |
| $y = D_{m,n}x$ | ```for (i=0;i<m*n;i++)```<br>```   y[i] = Dmn[i]*x[i];``` |
| $y = L_m^{mn}x$ | ```for (i=0;i<m;i++)```<br>```   for (j=0;j<n;j++)```<br>```      y[i+m*j]=x[n*i+j];``` |

**Example: tensor product**

$$\mathbf{I}_m \otimes A_n = \begin{bmatrix} A_n & & \\ & \ddots & \\ & & A_n \end{bmatrix}$$

# Program Generation in Spiral (Sketched)

**Transform**
**user specified**

$\mathsf{DFT}_8$

*Optimization at all abstraction levels*

**Fast algorithm**
**in SPL**
**many choices**

$$(\mathsf{DFT}_2 \otimes I_4)\, T_4^8 \left(I_2 \otimes \left((\mathsf{DFT}_2 \otimes I_2)\right.\right.$$
$$\left.\left. \cdot\, T_2^4\, (I_2 \otimes \mathsf{DFT}_2)\, L_2^4\right)\right) L_2^8$$

**parallelization**
**vectorization**

**∑-SPL:**

$$\sum \left(S_j \,\mathsf{DFT}_2\, G_j\right) \sum \left(\sum \left(S_{k,l}\,\mathrm{diag}\!\left(\mathsf{t}_{k,l}\right) \mathsf{DFT}_2\, G_l\right)\right.$$
$$\left.\sum \left(S_m\,\mathrm{diag}\!\left(\mathsf{t}_m\right) \mathsf{DFT}_2\, G_{k,m}\right)\right)$$

**loop**
**optimizations**

**C Code:**

```
void sub(double *y, double *x) {
  double f0, f1, f2, f3, f4, f7, f8, f10, f11;
  f0 = x[0] - x[3];
  f1 = x[0] + x[3];
  f2 = x[1] - x[2];
  f3 = x[1] + x[2];
  f4 = f1 - f3;
  y[0] = f1 + f3;
  y[2] = 0.7071067811865476 * f4;
  f7 = 0.9238795325112867 * f0;
  f8 = 0.3826834323650898 * f2;
  y[1] = f7 + f8;
  f10 = 0.3826834323650898 * f0;
  f11 = (-0.9238795325112867) * f2;
  y[3] = f10 + f11;
}
```

**constant folding**
**scheduling**
**……**

# Spiral Web Interface @spiral.net

**Program Generation Interface** collapse

**Target platform for optimization:** [2x Intel Xeon 3.6 GHz with 2048K L2 cache ▼]

**1. Select platform**

| parameter | value | explanation |
|---|---|---|
| **Transform** | DCT2 (Discrete Cosine Transform, type 2) ▼ | The transform for which you want to request C code |
| **Data type** | double ▼ | The data type of input and output vector |
| **Transform size** | 6 ▼ | The size of the transform = the length of the input vector |
| **Optimize for** | runtime ▼ | What you want to optimize the code for |
| **Search method** | Dynamic Programming ▼ | The search method SPIRAL uses (Dynamic Programming is a good choice) |
| **Compiler profile** | gcc -O3 ▼ | Compiler and compiler options used for compilation |

**2. Select functionality**

**3.** [ Generate code ] **"click"**

[ Generate code ]

**Browse Archive** expand

**Filter by Platform:** [All Platforms Selected ▼]

**Filter by Transform:** [All Transforms Selected ▼]

**Filter by Size:** [All Sizes Selected ▼]

[ Query Database ]

# Organization

- **Spiral overview**

- **Parallelization in Spiral**

- **Beyond Transforms**

- **Generating general-size libraries**

- **Results**

- **Concluding remarks**

# Types of Parallelism

- **Shared Memory (Multicore)**

- **Vector SIMD (SSE, VMX, Double FPU...)**

- **Message Passing (Clusters)**

- **Graphics Processors (GPUs)**

- **FPGA**

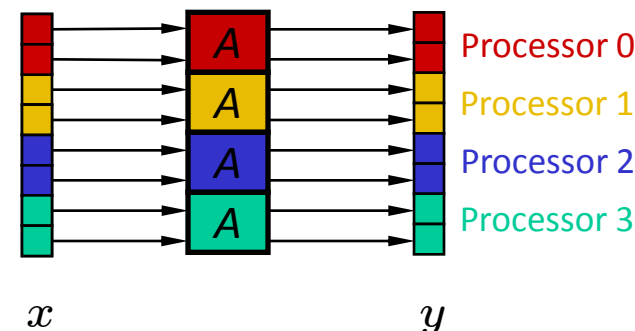- **HW/SW partitioning**

- **Multiple Levels of Parallelism (Cell BE)**

DSP transform (user specified)

Algorithm Generation
Algorithm Optimization

Formulas

Implementation
Code Optimization

Code

Compilation
Performance Evaluation

optimized/adapted implementation

*Spiral: One methodology optimizes for all types of parallelism*

# SPL to Shared Memory Code: Basic Idea

- **Governing construct: tensor product**

$$y = \left( \mathrm{I}_p \otimes A \right) x$$



Processor 0
Processor 1
Processor 2
Processor 3

$x$      $y$

**Independent operation, load-balanced**

- **Problematic construct: permutations produce false sharing**

$$y = \mathrm{L}_4^8 \, x$$



$x$      $y$

*Task: Rewrite formulas to*
*extract tensor product + keep contiguous blocks*

# Step 1: Shared Memory Tags

- **Identify crucial hardware parameters**
  - Number of processors: p
  - Cache line size: μ

- **Introduce them as tags in SPL**

$$\overbrace{A}\\ \mathsf{smp}(p,\mu)$$

**This means:** formula A is to be optimized for p processors and cache line size μ

# Step 2: Identify "Good" Formulas

■ **Load balanced, avoiding false sharing**

$$y = \Big( \mathrm{I}_p \otimes A \Big) x \quad \text{with} \quad A \in \mathbb{C}^{m\mu \times m\mu}$$

$$y = \left( \bigoplus_{i=0}^{p-1} A_i \right) x \quad \text{with} \quad A_i \in \mathbb{C}^{m\mu \times m\mu}$$

$$y = \Big( P \otimes \mathrm{I}_\mu \Big) x \quad \text{with } P \text{ a permutation matrix}$$

■ **Tagged operators (no further rewriting necessary)**

$$\mathrm{I}_p \otimes_\| A, \quad \bigoplus_{i=0}^{p-1} {}_\| A_i, \quad P \overline{\otimes} \mathrm{I}_\mu$$

■ **Definition: A formula is fully optimized if it is one of the above or of the form**

$$\mathrm{I}_m \otimes A \quad \text{or} \quad AB$$

**where A and B are fully optimized.**

# Step 3: Identify Rewriting Rules

- **Goal: Transform formulas into fully optimized formulas**
  - Formulas rewritten, tags propagated
  - There may be choices

$$\underbrace{AB}_{\mathsf{smp}(p,\mu)} \;\rightarrow\; \underbrace{A}_{\mathsf{smp}(p,\mu)}\;\underbrace{B}_{\mathsf{smp}(p,\mu)}$$

$$\underbrace{A_m \otimes \mathrm{I}_n}_{\mathsf{smp}(p,\mu)} \;\rightarrow\; \underbrace{\left(\mathsf{L}_m^{mp} \otimes \mathrm{I}_{n/p}\right)\left(\mathrm{I}_p \otimes (A_m \otimes \mathrm{I}_{n/p})\right)\left(\mathsf{L}_p^{mp} \otimes \mathrm{I}_{n/p}\right)}_{\mathsf{smp}(p,\mu)}$$

$$\underbrace{\mathsf{L}_m^{mn}}_{\mathsf{smp}(p,\mu)} \;\rightarrow\; \begin{cases} \underbrace{\left(\mathrm{I}_p \otimes \mathsf{L}_{m/p}^{mn/p}\right)}_{\mathsf{smp}(p,\mu)}\underbrace{\left(\mathsf{L}_p^{pn} \otimes \mathrm{I}_{m/p}\right)}_{\mathsf{smp}(p,\mu)} \\[2em] \underbrace{\left(\mathsf{L}_m^{pm} \otimes \mathrm{I}_{n/p}\right)}_{\mathsf{smp}(p,\mu)}\underbrace{\left(\mathrm{I}_p \otimes \mathsf{L}_m^{mn/p}\right)}_{\mathsf{smp}(p,\mu)} \end{cases}$$

$$\underbrace{\mathrm{I}_m \otimes A_n}_{\mathsf{smp}(p,\mu)} \;\rightarrow\; \mathrm{I}_p \otimes_{\|}\left(\mathrm{I}_{m/p} \otimes A_n\right)$$

$$\underbrace{(P \otimes \mathrm{I}_n)}_{\mathsf{smp}(p,\mu)} \;\rightarrow\; \left(P \otimes \mathrm{I}_{n/\mu}\right)\overline{\otimes}\,\mathrm{I}_\mu$$

# Parallelization by Rewriting

$$\underbrace{\mathbf{DFT}_{mn}}_{\mathsf{smp}(p,\mu)} \quad \rightarrow \quad \underbrace{\left((\mathbf{DFT}_m \otimes \mathbf{I}_n) \mathsf{T}_n^{mn} (\mathbf{I}_m \otimes \mathbf{DFT}_n) \mathsf{L}_m^{mn}\right)}_{\mathsf{smp}(p,\mu)}$$

$$\cdots$$

$$\rightarrow \quad \underbrace{\left(\mathbf{DFT}_m \otimes \mathbf{I}_n\right)}_{\mathsf{smp}(p,\mu)} \underbrace{\mathsf{T}_n^{mn}}_{\mathsf{smp}(p,\mu)} \underbrace{\left(\mathbf{I}_m \otimes \mathbf{DFT}_n\right)}_{\mathsf{smp}(p,\mu)} \underbrace{\mathsf{L}_m^{nm}}_{\mathsf{smp}(p,\mu)}$$

$$\cdots$$

$$\rightarrow \quad \left((\mathsf{L}_m^{mp} \otimes \mathbf{I}_{n/p\mu}) \otimes_\mu \mathbf{I}_\mu\right)\left(\mathbf{I}_p \otimes_\| (\mathbf{DFT}_m \otimes \mathbf{I}_{n/p})\right)\left((\mathsf{L}_p^{mp} \otimes \mathbf{I}_{n/p\mu}) \otimes_\mu \mathbf{I}_\mu\right)$$

$$\left(\bigoplus_{i=0}^{p-1}{}_\| \mathsf{T}_n^{mn,i}\right)\left(\mathbf{I}_p \otimes_\| (\mathbf{I}_{m/p} \otimes \mathbf{DFT}_n)\right)\left(\mathbf{I}_p \otimes_\| \mathsf{L}_{m/p}^{mn/p}\right)\left((\mathsf{L}_p^{pn} \otimes \mathbf{I}_{m/p\mu}) \otimes_\mu \mathbf{I}_\mu\right)$$

**Fully optimized (load-balanced, no false sharing)**
**in the sense of our definition**

# Same Approach for Other Parallel Paradigms

## Threading:

$$\underline{\frac{DFT_{mn}}{smp(p,\mu)}} \rightarrow \underline{\left((DFT_m \otimes I_n)T_n^{mn}(I_m \otimes DFT_n)L_m^{mn}\right)}_{smp(p,\mu)}$$

$$\dots$$

$$\rightarrow \underline{\left(DFT_m \otimes I_n\right)}_{smp(p,\mu)} \underline{T_n^{mn}}_{smp(p,\mu)} \underline{\left(I_m \otimes DFT_n\right)}_{smp(p,\mu)} \underline{L_m^{nm}}_{smp(p,\mu)}$$

$$\dots$$

$$\rightarrow \left((L_m^{mp} \otimes I_{n/p\mu}) \otimes_\mu I_\mu\right)\left(I_p \otimes_\parallel (DFT_m \otimes I_{n/p})\right)\left((L_p^{mp} \otimes I_{n/p\mu}) \otimes_\mu I_\mu\right)$$

$$\left(\bigoplus_{i=0}^{p-1} {}_\parallel T_n^{mn,i}\right)\left(I_p \otimes_\parallel (I_{m/p} \otimes DFT_n)\right)\left(I_p \otimes_\parallel L_{m/p}^{mn/p}\right)\left((L_p^{pn} \otimes I_{m/p\mu}) \otimes_\mu I_\mu\right)$$

## Vectorization:

$$\underline{\overline{\left(\frac{DFT_{mn}}{vec(\nu)}\right)}} \rightarrow \underline{\left((DFT_m \otimes I_n)T_n^{mn}(I_m \otimes DFT_n)L_m^{mn}\right)}_{vec(\nu)}$$

$$\dots$$

$$\rightarrow \overleftarrow{\underline{(DFT_m \otimes I_n)}_{vec(\nu)}}^\nu \overleftarrow{\underline{(T_n^{mn})}_{vec(\nu)}}^\nu \overrightarrow{\underline{(I_m \otimes DFT_n)L_m^{mn}}_{vec(\nu)}}^\nu$$

$$\dots$$

$$\rightarrow (I_{mn/\nu} \otimes \underline{L_\nu^{2\nu}}_{sse})(\overline{DFT_m \otimes I_{n/\nu}} \vec\otimes I_\nu) \overline{(T_n^{mn})}_{sse}^\nu$$

$$\left(I_{m/\nu} \otimes (\overline{L_\nu^n} \vec\otimes I_\nu)(I_{n/\nu} \otimes (L_\nu^{2\nu} \vec\otimes I_\nu)(I_2 \otimes \underline{L_\nu^{\nu^2}}_{sse})(L_2^{2\nu} \vec\otimes I_\nu))(\overline{DFT_n} \vec\otimes I_\nu)\right)$$

$$\left((L_m^{mn} \otimes I_2) \vec\otimes I_\nu\right)(I_{mn/\nu} \otimes \underline{L_2^{2\nu}}_{sse})$$

## GPUs:

$$\underline{\left(\frac{DFT_{r^k}}{gpu(t,c)}\right)} \rightarrow \underline{\left(\prod_{i=0}^{k-1} L_r^{r^k}\left(I_{r^{k-1}} \otimes DFT_r\right)\left(L_{r^{k-i-1}}^{r^k}(I_{r^i} \otimes T_{r^{k-i-1}}^{r^{k-i}})\underline{L_{r^{i+1}}^{r^k}}_{vec(c)}\right)\right)R_r^{r^k}}_{gpu(t,c)}$$

$$\dots$$

$$\rightarrow \left(\prod_{i=0}^{k-1}(L_r^{r^n/2} \vec\otimes I_2)\left(I_{r^{n-1}/2} \otimes_\times \underline{(DFT_r \vec\otimes I_2)L_r^{2r}}_{shd(t,c)}\right)T_i\right)$$

$$(L_r^{r^n/2} \vec\otimes I_2)(I_{r^{n-1}/2} \otimes_\times \underline{L_r^{2r}}_{shd(t,c)})(R_r^{r^{n-1}} \vec\otimes I_r)$$

## Verilog for FPGAs:

$$\underline{\left(\frac{DFT_{r^k}}{stream(r^s)}\right)} \rightarrow \underline{\left[\prod_{i=0}^{k-1} L_r^{r^k}\left(I_{r^{k-1}} \otimes DFT_r\right)\left(L_{r^{k-i-1}}^{r^k}(I_{r^i} \otimes T_{r^{k-i-1}}^{r^{k-i}})L_{r^{i+1}}^{r^k}\right)\right]R_r^{r^k}}_{stream(r^s)}$$

$$\dots$$

$$\rightarrow \left[\prod_{i=0}^{k-1} \underline{L_r^{r^k}}_{stream(r^s)} \underline{\frac{\left(I_{r^{k-1}} \otimes DFT_r\right)}{stream(r^s)}} \underline{\frac{\left(L_{r^{k-i-1}}^{r^k}(I_{r^i} \otimes T_{r^{k-i-1}}^{r^{k-i}})L_{r^{i+1}}^{r^k}\right)}{stream(r^s)}}\right]\underline{R_r^{r^k}}_{stream(r^s)}$$

$$\dots$$

$$\rightarrow \left[\prod_{i=0}^{k-1} \underline{L_r^{r^k}}_{stream(r^s)} \underline{\left(I_{r^{k-s-1}} \otimes_s (I_{r^{s-1}} \otimes DFT_r)\right)} \underline{T_i'}_{stream(r^s)}\right]\underline{R_r^{r^k}}_{stream(r^s)}$$

■ **Rigorous, correct by construction**

■ **Overcomes compiler limitations**

SPIRAL
www.spiral.net

# Pre-Silicon Optimization: Larrabee and AVX

```c
void dft64(float  *Y, float  *X) {
    __m512 U912, U913, U914, U915, U916, U917, U918, U919, U920, U921, U922, U923, U924, U925,...
 __m512  *a2153, *a2155;
    a2153 = ((__m512  *) X);   s1107 = *(a2153);
    s1108 = *((a2153 + 4));    t1323 = _mm512_add_ps(s1107,s1108);
    t1324 = _mm512_sub_ps(s1107,s1108);
    ...
    U926 = _mm512_swizupconv_r32(_mm512_set_1to16_ps(0.70710678118654757),_MM_SWIZ_REG_CDAB);
    s1121 = _mm512_madd231_ps(_mm512_mul_ps(_mm512_mask_or_pi(
        _mm512_set_1to16_ps(0.70710678118654757),0xAAAA,a2154,U926),t1341),
        _mm512_mask_sub_ps(_mm512_set_1to16_ps(0.70710678118654757),0x5555,a2154,U926),
        _mm512_swizupconv_r32(t1341,_MM_SWIZ_REG_CDAB));
    U927 = _mm512_swizupconv_r32(_mm512_set_16to16_ps(0.70710678118654757, (-0.70710678118654757),
        0.70710678118654757, (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757),
        0.70710678118654757, (-0.70710678118654757), 0.70710678118654757, (-0.70710678118654757),
        0.70710678118654757, (-0.70710678118654757
        0.70710678118654757, (-0.70710678118654757
    ...
    s1166 = _mm512_madd231_ps(_mm512_mul_ps(_mm51
        0.70710678118654757, (-0.70710678118654757
        0.70710678118654757, (-0.70710678118654757
        0.70710678118654757, (-0.70710678118654757
        0.70710678118654757, (-0.70710678118654757
        0xAAAA,a2154,U951),t1362),
        _mm512_mask_sub_ps(_mm512_set_16to16_ps(0
        (-0.70710678118654757), 0.707106781186547
        (-0.70710678118654757), 0.707106781186547
        (-0.70710678118654757), 0.707106781186547
        (-0.70710678118654757), 0.707106781186547
        _mm512_swizupconv_r32(t1362,_MM_SWIZ_REG_
    ...
}
```

**DFT (double-precision), single core Larrabee Native**
performance [Gflop/s]



DFT on Larrabee

input size

**Not actual data (NDA)**

# Organization

- **Spiral overview**

- **Parallelization in Spiral**

- **Beyond Transforms**

- **Generating general-size libraries**

- **Results**

- **Concluding remarks**

F. Franchetti, F. de Mesmay, D. McFarlin, and M. Püschel:
**Operator Language: A Program Generation Framework for Fast Kernels.** In Proceedings of DSL WC, 2009.

# Going Beyond Transforms

- **Transform =
  linear operator with one vector input and one vector output**



- **Key ideas:**
  - Generalize to (possibly nonlinear) operators with several inputs and several outputs
  - Generalize SPL (including tensor product) to OL (operator language)
  - Generalize rewriting systems for parallelizations

# Operator Language

| name | definition |
|------|------------|

*basic operators*

| | |
|---|---|
| projection | $\pi_{\mathbf{x}} : \mathbb{C}^m \times \mathbb{C}^n \to \mathbb{C}^m; \ (\mathbf{x}, \mathbf{y}) \mapsto \mathbf{x}$ |
| linear transform | $\mathsf{M} : \mathbb{C}^n \to \mathbb{C}^m; \mathbf{x} \mapsto M\mathbf{x}$ |
| stride | $\mathsf{L}_m^{mn} : \mathbb{C}^{mn} \to \mathbb{C}^{mn}; \ \mathbf{x} \mapsto \mathsf{L}_m^{mn}\mathbf{x}$ |
| vector sum | $\Sigma_n : \mathbb{C}^n \to \mathbb{C}; \ \mathbf{x} \mapsto \sum_{i=0}^{n-1} x_i$ |
| vector minimum | $\min_n : \mathbb{C}^n \to \mathbb{C}; \ \mathbf{x} \mapsto \min(x_0, \ldots, x_{n-1})$ |
| constant vector | $\mathsf{C}_{\mathbf{c}} : \varnothing \to \mathbb{C}^n; \ () \mapsto \mathbf{c}$ |

*operations*

| | |
|---|---|
| addition | $(\mathsf{M} + \mathsf{N})(\mathbf{x}, \mathbf{y}) = \mathsf{M}(\mathbf{x}, \mathbf{y}) + \mathsf{N}(\mathbf{x}, \mathbf{y})$ |
| multiplication | $(\mathsf{M} \cdot \mathsf{N})(\mathbf{x}, \mathbf{y}) = \mathsf{M}(\mathbf{x}, \mathbf{y}) \cdot \mathsf{N}(\mathbf{x}, \mathbf{y})$ |
| direct sum | $(\mathsf{M} \oplus \mathsf{N})(\mathbf{x} \oplus \mathbf{u}, \mathbf{y} \oplus \mathbf{v}) = \mathsf{M}(\mathbf{x}, \mathbf{y}) \oplus \mathsf{N}(\mathbf{u}, \mathbf{v})$ |
| cartesian product | $(\mathsf{M} \times \mathsf{N})(\mathbf{x}, \mathbf{y}, \mathbf{u}, \mathbf{v}) = \mathsf{M}(\mathbf{x}, \mathbf{y}) \times \mathsf{N}(\mathbf{u}, \mathbf{v})$ |
| composition | $(\mathsf{M} \circ \mathsf{N})(\mathbf{x}, \mathbf{y}) = \mathsf{M}(\mathsf{N}(\mathbf{x}, \mathbf{y}))$ |
| iterative composition | $\left( \prod_{i=0}^{n-1} \mathsf{M}_i \right)(\mathbf{x}, \mathbf{y}) = (\mathsf{M}_0 \circ \cdots \circ \mathsf{M}_{n-1})(\mathbf{x}, \mathbf{y})$ |
| tensor product | $\mathsf{I} \otimes \mathsf{M}, \ \mathsf{M} \otimes \mathsf{I}$ |

$$(\mathsf{I}_{m \times n \to mn} \otimes \mathsf{A}) \left( \sum_{i=0}^{m-1} \mathbf{e}_i^m \otimes \mathbf{x}, \sum_{i=0}^{n-1} \mathbf{e}_i^n \otimes \mathbf{y} \right) = \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} \mathbf{e}_i^m \otimes \mathbf{e}_j^n \otimes \mathsf{A}(\mathbf{x}, \mathbf{y}),$$

$$(\mathsf{A} \otimes \mathsf{I}_{m \times n \to mn}) \left( \sum_{i=0}^{m-1} \mathbf{x} \otimes \mathbf{e}_i^m, \sum_{i=0}^{n-1} \mathbf{y} \otimes \mathbf{e}_i^n \right) = \sum_{i=0}^{m-1} \sum_{j=0}^{n-1} \mathsf{A}(\mathbf{x}, \mathbf{y}) \otimes \mathbf{e}_i^m \otimes \mathbf{e}_j^n.$$

# Expressing Kernels as OL Formulas

## Linear Transforms

$$\mathbf{DFT}_n \;\to\; (\mathbf{DFT}_k \otimes \mathrm{I}_m)\, \mathsf{T}^n_m (\mathrm{I}_k \otimes \mathbf{DFT}_m)\, \mathsf{L}^n_k, \quad n = km$$

$$\mathbf{DFT}_n \;\to\; P_n(\mathbf{DFT}_k \otimes \mathbf{DFT}_m)Q_n, \quad n = km,\ \gcd(k,m)=1$$

$$\mathbf{DFT}_p \;\to\; R_p^T(\mathrm{I}_1 \oplus \mathbf{DFT}_{p-1})D_p(\mathrm{I}_1 \oplus \mathbf{DFT}_{p-1})R_p, \quad p \text{ prime}$$

$$\mathbf{DCT\text{-}3}_n \;\to\; (\mathrm{I}_m \oplus \mathsf{J}_m)\, \mathsf{L}^n_m (\mathbf{DCT\text{-}3}_m(1/4) \oplus \mathbf{DCT\text{-}3}_m(3/4))$$

$$\cdot (\mathsf{F}_2 \otimes \mathrm{I}_m) \begin{bmatrix} \mathrm{I}_m & 0 \oplus -\mathsf{J}_{m-1} \\ \frac{1}{\sqrt{2}}(\mathrm{I}_1 \oplus 2\,\mathrm{I}_m) \end{bmatrix}, \quad n = 2m$$

$$\mathbf{DCT\text{-}4}_n \;\to\; S_n\mathbf{DCT\text{-}2}_n \operatorname{diag}_{0\le k<n}(1/(2\cos((2k+1)\pi/4n)))$$

$$\mathbf{IMDCT}_{2m} \;\to\; (\mathsf{J}_m \oplus \mathrm{I}_m \oplus \mathrm{I}_m \oplus \mathsf{J}_m)\left(\left(\begin{bmatrix}1\\-1\end{bmatrix}\otimes \mathrm{I}_m\right)\oplus\left(\begin{bmatrix}-1\\-1\end{bmatrix}\otimes \mathrm{I}_m\right)\right)\mathsf{J}_{2m}\mathbf{DCT\text{-}4}_{2m}$$

$$\mathbf{WHT}_{2^k} \;\to\; \prod_{i=1}^{t}(\mathrm{I}_{2^{k_1+\cdots+k_{i-1}}} \otimes \mathbf{WHT}_{2^{k_i}} \otimes \mathrm{I}_{2^{k_{i+1}+\cdots+k_t}}), \quad k = k_1 + \cdots + k_t$$

$$\mathbf{DFT}_2 \;\to\; \mathsf{F}_2$$

$$\mathbf{DCT\text{-}2}_2 \;\to\; \operatorname{diag}(1, 1/\sqrt{2})\, \mathsf{F}_2$$

$$\mathbf{DCT\text{-}4}_2 \;\to\; \mathsf{J}_2\, \mathsf{R}_{13\pi/8}$$

## Viterbi Decoding



$$\underbrace{\mathbf{Vit}}_{\mathrm{vec}(v)} \;\to\; \underbrace{\left(\prod (L \times I) \circ (I \otimes C)\right) \circ Id}_{\mathrm{vec}(v)}$$

$$\to\; \left(\prod \underbrace{(L \times I) \circ (I \otimes C)}_{\mathrm{vec}(v)}\right) \circ Id$$

$$\times \;\to\; \left(\prod (L \otimes I_v \times I) \circ (I \otimes C \otimes I_v) \circ (\vec{L} \times I)\right) \circ Id$$

$$\to\; \prod (L \otimes I_v \times I) \circ (I \otimes (B \otimes I_v)) \circ (\vec{L} \times I)$$

## Matrix-Matrix Multiplication



$$\mathbf{MMM}_{1,1,1} \to (\cdot)_1$$

$$\mathbf{MMM}_{m,n,k} \to (\otimes)_{m/m_b \times 1} \otimes \mathbf{MMM}_{m_b,n,k}$$

$$\mathbf{MMM}_{m,n,k} \to \mathbf{MMM}_{m,nb,k} \otimes (\otimes)_{1\times n/nb}$$

$$\mathbf{MMM}_{m,n,k} \to ((\Sigma_{k/k_b} \circ (\cdot)_{k/k_b}) \otimes \mathbf{MMM}_{m,n,k_b})\circ$$
$$((L^{mk/k_b}_{k/k_b} \otimes I_{k_b}) \times I_{kn})$$

$$\mathbf{MMM}_{m,n,k} \to (L^{mn/n_b}_m \otimes I_{n_b})\circ$$
$$((\otimes)_{1\times n/n_b} \otimes \mathbf{MMM}_{m,n_b,k})\circ$$
$$(I_{km} \times (L^{kn/n_b}_{n/n_b} \otimes I_{n_b}))$$

## Synthetic Aperture Radar (SAR)



$$\mathbf{SAR}_{k\times m \to n\times n} \;\to\; \mathbf{DFT}_{n\times n} \circ \mathbf{Interp}_{k\times m \to n\times n}$$

$$\mathbf{DFT}_{n\times n} \;\to\; (\mathbf{DFT}_n \otimes \mathrm{I}_n) \circ (\mathrm{I}_n \otimes \mathbf{DFT}_n)$$

$$\mathbf{Interp}_{k\times m \to n\times n} \;\to\; (\mathbf{Interp}_{k\to n} \otimes_i \mathrm{I}_n) \circ (\mathrm{I}_k \otimes_i \mathbf{Interp}_{m\to n})$$

$$\mathbf{Interp}_{r\to s} \;\to\; \left(\bigoplus_{i=0}^{n-2} \mathbf{InterpSeg}_k\right) \oplus \mathbf{InterpSegPruned}_{k,\ell}$$

$$\mathbf{InterpSeg}_k \;\to\; \mathsf{G}^{u\cdot n\to k}_f \circ \mathbf{iPrunedDFT}_{n\to u\cdot n} \circ \left(\frac{1}{n}\right) \circ \mathbf{DFT}_n$$

# Example: Matrix Multiplication (MMM)

**Breakdown rules:**

*capture various forms of blocking*

| breakdown rule | description |
|---|---|
| $\text{MMM}_{1,1,1} \rightarrow (\cdot)_1$ | base case |
| $\text{MMM}_{m,n,k} \rightarrow (\otimes)_{m/m_b \times 1} \otimes \text{MMM}_{m_b,n,k}$ | horizontal blocking |
| $\text{MMM}_{m,n,k} \rightarrow \text{MMM}_{m,nb,k} \otimes (\otimes)_{1 \times n/nb}$ | interleaved blocking |
| $\text{MMM}_{m,n,k} \rightarrow ((\Sigma_{k/k_b} \circ (\cdot)_{k/k_b}) \otimes \text{MMM}_{m,n,k_b}) \circ$ $\qquad ((L_{k/k_b}^{mk/k_b} \otimes I_{k_b}) \times I_{kn})$ | accumulative blocking |
| $\text{MMM}_{m,n,k} \rightarrow (L_m^{mn/n_b} \otimes I_{n_b}) \circ$ $\qquad ((\otimes)_{1 \times n/n_b} \otimes \text{MMM}_{m,n_b,k}) \circ$ $\qquad (I_{km} \times (L_{n/n_b}^{kn/n_b} \otimes I_{n_b}))$ | vertical blocking |

# Parallelization: OL Rewriting Rules

- **SPL rules and hardware model reused**
- **Few additional OL-specific rules required**

$$\underbrace{\left( I_k \otimes L_n^{mn} \right)}_{\mathsf{smp}(p,\mu)} \circ \underbrace{L_{km}^{kmn}}_{\mathsf{smp}(p,\mu)} \rightarrow \left( L_k^{kn} \otimes I_{m/\mu} \right) \bar{\otimes} I_\mu$$

$$\underbrace{L_n^{kmn}}_{\mathsf{smp}(p,\mu)} \circ \underbrace{\left( I_k \otimes L_m^{mn} \right)}_{\mathsf{smp}(p,\mu)} \rightarrow \left( L_n^{kn} \otimes I_{m/\mu} \right) \bar{\otimes} I_\mu$$

$$\underbrace{A \circ B}_{\mathsf{smp}(p,\mu)} \rightarrow \underbrace{A}_{\mathsf{smp}(p,\mu)} \circ \underbrace{B}_{\mathsf{smp}(p,\mu)}$$

$$\underbrace{A^{k \times m \to n} \otimes I^{1 \times p \to p}}_{\mathsf{smp}(p,\mu)} \rightarrow \underbrace{L_n^{pn}}_{\mathsf{smp}(p,\mu)} \circ \left( I_{1 \times p \to p} \otimes_{\|} A^{k \times m \to n} \right) \circ \underbrace{\left( I_k \times L_p^{pm} \right)}_{\mathsf{smp}(p,\mu)}$$

$$\underbrace{\left( A \times B \right)}_{\mathsf{smp}(p,\mu)} \circ \underbrace{\left( C \times D \right)}_{\mathsf{smp}(p,\mu)} \rightarrow \underbrace{\left( A \circ C \right)}_{\mathsf{smp}(p,\mu)} \times \underbrace{\left( B \circ D \right)}_{\mathsf{smp}(p,\mu)}$$

**New OL rules**

# Parallelization Through Rewriting: MMM

$$\underbrace{\mathrm{MMM}_{m,n,k}}_{\mathrm{smp}(p,\mu)}$$

$$\rightarrow \underbrace{\left(\mathrm{I}_m \otimes \mathsf{L}_p^n\right) \circ \left(\mathrm{MMM}_{m,n/p,k} \otimes (\otimes)_{1\times p \rightarrow p}\right) \circ \left(\mathrm{I}_{km} \times (\mathrm{I}_k \otimes \mathsf{L}_{n/p}^n)\right)}_{\mathrm{smp}(p,\mu)}$$

$$\rightarrow \underbrace{\left(\mathrm{I}_m \otimes \mathsf{L}_p^n\right)}_{\mathrm{smp}(p,\mu)} \circ \underbrace{\left(\mathrm{MMM}_{m,n/p,k} \otimes (\otimes)_{1\times p \rightarrow p}\right)}_{\mathrm{smp}(p,\mu)} \circ \underbrace{\left(\mathrm{I}_{km} \times (\mathrm{I}_k \otimes \mathsf{L}_{n/p}^n)\right)}_{\mathrm{smp}(p,\mu)}$$

$$\rightarrow \underbrace{\left(\mathrm{I}_m \otimes \mathsf{L}_p^n\right)}_{\mathrm{smp}(p,\mu)} \circ \underbrace{\mathsf{L}_{m/pn}^{mn}}_{\mathrm{smp}(p,\mu)} \circ \left((\otimes)_{1\times p \rightarrow p} \otimes_\| \mathrm{MMM}_{m/p,n,k}\right) \circ \underbrace{\left(\mathrm{I}_{km} \times \mathsf{L}_p^{kn}\right)}_{\mathrm{smp}(p,\mu)} \circ \underbrace{\left(\mathrm{I}_{km} \times (\mathrm{I}_k \otimes \mathsf{L}_{n/p}^n)\right)}_{\mathrm{smp}(p,\mu)}$$

$$\rightarrow \left((\mathsf{L}_m^{mp} \otimes \mathrm{I}_{n/(p\mu)}) \bar\otimes \mathrm{I}_\mu\right) \circ \left((\otimes)_{1\times p \rightarrow p} \otimes_\| \mathrm{MMM}_{m,n/p,k}\right) \circ \left((\mathrm{I}_{km/\mu} \bar\otimes \mathrm{I}_\mu) \times ((\mathsf{L}_p^{kp} \otimes \mathrm{I}_{n/(p\mu)}) \bar\otimes \mathrm{I}_\mu)\right)$$

**Load-balanced**
**No false sharing**

# Organization

■ **Spiral overview**

■ **Parallelization in Spiral**

■ **Beyond Transforms**

■ **Generating general-size libraries**
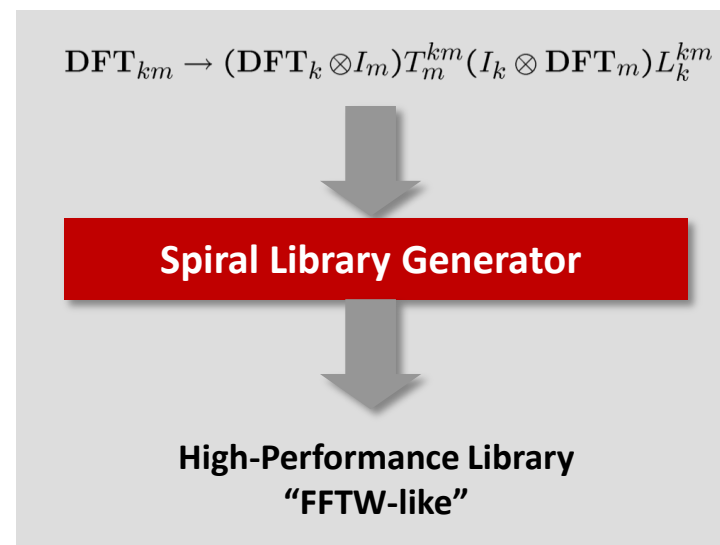
■ **Results**

■ **Concluding remarks**

Y. Voronenko, F. de Mesmay, and M. Püschel: **Computer generation of general size linear transform libraries**
in Proceedings Code Generation and Optimization (CGO), 2009.

# General-Size Library

## Input:

- **Transform**: $\mathrm{DFT}_n$

- **Algorithms**: $\mathrm{DFT}_{km} \to (\mathrm{DFT}_k \otimes I_m) T_m^{km} (I_k \otimes \mathrm{DFT}_m) L_k^{km}$
  $\mathrm{DFT}_2 \to \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$

- **Vectorization**: 2-way SSE

- **Threading**: Yes

## Output:

- Optimized library (10,000 lines of C++)

- For general input size
  (**not** collection of fixed sizes)

- Vectorized

- Multithreaded

- With runtime adaptation mechanism

- Performance competitive with hand-written code

$\mathrm{DFT}_{km} \to (\mathrm{DFT}_k \otimes I_m) T_m^{km} (I_k \otimes \mathrm{DFT}_m) L_k^{km}$

**Spiral Library Generator**

**High-Performance Library
"FFTW-like"**

# Beyond Fourier Transform and FFTW

| Cooley-Tukey FFT | "Cooley-Tukey" DCT | Overlap-save/add FIR |
|:---:|:---:|:---:|
| ↓ | ↓ | ↓ |
| **Spiral** | **Spiral** | **Spiral** |
| ↓ | ↓ | ↓ |
| "FFTW" | "FCTW" | "FIRW" |

| Fast Walsh Transform | Fast Wavelet Transform | Fast Hartley Transform |
|:---:|:---:|:---:|
| ↓ | ↓ | ↓ |
| **Spiral** | **Spiral** | **Spiral** |
| ↓ | ↓ | ↓ |
| "WHTW" | "FWTW" | "FHTW" |

**Y. Voronenko's PhD Thesis (experimental results): 50+ "FFTW-like" libraries**

# How it Works: Recursion Step Closure

■ **Input:** transform T and a breakdown rule

■ **Output:** recursion steps + Σ-SPL implementation

■ **Algorithm:**

$$\{\mathbf{DFT}_n\}$$

1. Apply the breakdown rule

$$(\{\mathbf{DFT}_{n/k}\} \otimes I_k)T_k^n(I_{n/k} \otimes \{\mathbf{DFT}_k\})L_{n/k}^n$$

2. Convert to Σ-SPL

$$\left(\sum_{i=0}^{k-1} \mathsf{S}_{h_{i,k}}\{\mathbf{DFT}_{n/k}\}\, \mathsf{G}_{h_{i,k}}\right) \mathrm{diag}(f) \left(\sum_{j=0}^{n/k-1} \mathsf{S}_{h_{jk,1}}\{\mathbf{DFT}_k\}\, \mathsf{G}_{h_{jk,1}}\right)\mathrm{perm}(\ell_{n/k}^n)$$

3. Apply loop merging + index simplification rules.

$$\sum_{i=0}^{k-1} \mathsf{S}_{h_{i,k}}\{\mathbf{DFT}_{n/k}\}\,\mathrm{diag}(f\circ h_{i,k})\, \mathsf{G}_{h_{i,k}} \sum_{j=0}^{n/k-1} \mathsf{S}_{h_{jk,1}}\{\mathbf{DFT}_k\}\, \mathsf{G}_{h_{j,n/k}}$$

4. Extract recursion steps

$$\sum_{i=0}^{k-1} \left\{ \mathsf{S}_{h_{i,k}}\, \mathbf{DFT}_{n/k}\,\mathrm{diag}(f\circ h_{i,k})\, \mathsf{G}_{h_{i,k}} \right\} \sum_{j=0}^{n/k-1} \left\{ \mathsf{S}_{h_{jk,1}}\, \mathbf{DFT}_k\, \mathsf{G}_{h_{j,n/k}} \right\}$$

5. Repeat until closure is reached

**Spiral translates recursion steps into adaptive library (codelets + recursion)**

# Recursion Step Closure: Examples

**DFT (scalar)**



**4 mutually recursive functions**
- **computed automatically**
- **described using Σ-SPL formulas**

**DCT4 (vectorized)**



$1:\ \mathrm{Vec}_2(\mathbf{DCT\text{-}4}_{u_1})$

$2:\ \mathrm{Vec}_2(\mathrm{GT}(\mathrm{diag}\,(N_{2u_8})\,\mathbf{RDFT\text{-}3}_{2u_8}^\top\,\mathrm{rcdiag}(\mathrm{pre}(u_4^{\mathbb{Z}\times 2u_8\to\mathbb{R}})),\ h_{0,\,1,u_7}^{2u_8\to u_6}\circ\ell_{u_8}^{2u_8},\ r_{0,\,u_{11},1,\,u_{12}}^{2u_8\to u_9},\ \{u_{13}\}))$

$3:\ \mathrm{Vec}_2(\mathrm{GT}(\mathbf{RDFT\text{-}3}_{u_1}\,\mathrm{diag}\,(N_{u_1}),\ r_{0,\,u_5,1,\,u_6}^{u_1\to u_3},\ h_{0,\,u_9,1}^{u_1\to u_8},\ \{u_{10}\}))$

$4:\ \mathrm{VJam}_2(\mathrm{GT}(\mathrm{diag}\,(N_{2u_9})\,\mathbf{RDFT\text{-}3}_{2u_9}^\top\,\mathrm{rcdiag}(\mathrm{pre}(u_4^{\mathbb{Z}\times\mathbb{Z}\times 2u_9\to\mathbb{R}})),\ h_{0,\,1,u_7,u_8}^{2u_9\to u_6}\circ\ell_{u_9}^{2u_9},\ r_{0,\,u_{12},1,2,\,u_{13}}^{2u_9\to u_{10}},\ \{2,u_{14}\}))$

$5:\ \mathrm{GT}(\mathrm{diag}\,(N_{2u_9})\,\mathbf{RDFT\text{-}3}_{2u_9}^\top\,\mathrm{rcdiag}(\mathrm{pre}(u_4^{\mathbb{Z}\times 2u_9\to\mathbb{R}})),\ h_{u_7,\,1,u_8}^{2u_9\to u_6}\circ\ell_{u_9}^{2u_9},\ r_{u_{12},\,u_{13},1,\,u_{14}}^{2u_9\to u_{10}},\ \{u_{15}\})$

$6:\ \mathrm{VJam}_2(\mathrm{GT}(\mathbf{RDFT\text{-}3}_{u_1}\,\mathrm{diag}\,(N_{u_1}),\ r_{0,\,u_5,1,2,\,u_6}^{u_1\to u_3},\ h_{0,\,u_9,1,2}^{u_1\to u_8},\ \{2,u_{10}\}))$

$7:\ \mathrm{GT}(\mathbf{RDFT\text{-}3}_{u_1}\,\mathrm{diag}\,(N_{u_1}),\ r_{u_5,\,u_6,1,\,u_7}^{u_1\to u_3},\ h_{u_{10},\,u_{11},1}^{u_1\to u_9},\ \{u_{12}\})$

$8:\ \mathrm{S}(h_{u_3,\,u_4}^{u_1\to u_2})\,\mathbf{RDFT\text{-}3}_{u_1}\,\mathrm{diag}\,(N_{u_1})\,\mathrm{G}(r_{u_9,\,u_{10},\,u_{11}}^{u_1\to u_7})$

$9:\ \mathrm{S}(r_{u_3,\,u_4,\,u_5}^{2u_{13}\to u_1})\mathrm{diag}\,(N_{2u_{13}})\,\mathbf{RDFT\text{-}3}_{2u_{13}}^\top\,\mathrm{rcdiag}(\mathrm{pre}(u_9^{2u_{13}\to\mathbb{R}}))\mathrm{G}(h_{u_{12},\,1}^{2u_{13}\to u_{11}}\circ\ell_{u_{13}}^{2u_{13}})$

$10:\ \mathrm{VJam}_2(\mathrm{GT}(\mathrm{diag}\,(N_{2u_9})\,\mathbf{RDFT\text{-}3}_{2u_9}^\top\,\mathrm{rcdiag}(\mathrm{pre}(u_4^{\mathbb{Z}\times 2u_9\to\mathbb{R}})),\ h_{u_7,\,1,u_8}^{2u_9\to u_6}\circ\ell_{u_9}^{2u_9},\ r_{u_{12},\,u_{13},1,\,u_{14}}^{2u_9\to u_{10}},\ \{2\}))$

$11:\ \mathrm{VJam}_2(\mathrm{GT}(\mathbf{RDFT\text{-}3}_{u_1}\,\mathrm{diag}\,(N_{u_1}),\ r_{u_5,\,u_6,1,\,u_7}^{u_1\to u_3},\ h_{u_{10},\,u_{11},1}^{u_1\to u_9},\ \{2\}))$

$12:\ \mathrm{GT}(\mathrm{diag}\,(C_{u_1})\,\mathbf{rDFT}_{2u_1}(\lambda\text{-}\mathrm{wrap}(\lambda_1^{\mathbb{Z}\times\mathbb{Z}\to\mathbb{R}})),\ h_{0,\,1,u_5}^{2u_1\to u_4},\ h_{u_8,\,u_9}^{2u_{10}\to u_7}\circ(r_{0,\,u_{12},1,\,u_{13}}^{u_1\to u_{10}}\otimes\imath_2),\ \{u_{14}\})$

$13:\ \mathrm{VJam}_2(\mathrm{GT}(\mathrm{diag}\,(C_{u_1})\,\mathbf{rDFT}_{2u_1}(\lambda\text{-}\mathrm{wrap}(\lambda_1^{\mathbb{Z}\times\mathbb{Z}\times\mathbb{Z}\to\mathbb{R}})),\ h_{u_5,\,u_6,1,u_7}^{2u_1\to u_4},\ h_{u_{10},\,u_{11},1}^{2u_{12}\to u_9}\circ(r_{0,\,u_{14},0,1,\,u_{15}}^{u_1\to u_{12}}\otimes\imath_2),\ \{2,u_{16}\}))$

$14:\ \mathrm{VJam}_2(\mathrm{GT}(\mathbf{RDFT\text{-}3}_{u_1}\,\mathrm{diag}\,(N_{u_1}),\ r_{u_5,\,u_6,1,u_7,\,u_8}^{u_1\to u_3},\ h_{u_{11},\,u_{12},1,u_{13}}^{u_1\to u_{10}},\ \{2,u_{14}\}))$

$15:\ \mathrm{GT}(\mathbf{RDFT\text{-}3}_{u_1}\,\mathrm{diag}\,(N_{u_1}),\ r_{u_5,\,u_6,u_7,\,u_8}^{u_1\to u_3},\ h_{0,\,u_{11},1}^{u_1\to u_{10}},\ \{u_{12}\})$

$16:\ \mathrm{S}(h_{u_3,\,u_4}^{2u_5\to u_2}\circ(r_{u_7,\,u_8,\,u_9}^{u_6\to u_5}\otimes\imath_2))\mathrm{diag}\,(C_{u_6})\,\mathbf{rDFT}_{2u_6}(\lambda\text{-}\mathrm{wrap}(\lambda_1^{\mathbb{Z}\to\mathbb{R}}))\mathrm{G}(h_{u_{14},\,1}^{2u_6\to u_{13}})$

$17:\ \mathrm{VJam}_2(\mathrm{GT}(\mathrm{diag}\,(C_{u_1})\,\mathbf{rDFT}_{2u_1}(\lambda\text{-}\mathrm{wrap}(\lambda_1^{\mathbb{Z}\times\mathbb{Z}\to\mathbb{R}})),\ h_{u_5,\,u_6,1}^{2u_1\to u_4},\ h_{u_9,\,u_{10},1}^{2u_{11}\to u_8}\circ(r_{u_{13},\,u_{14},\,u_{15}}^{u_1\to u_{11}}\otimes\imath_2),\ \{2\}))$
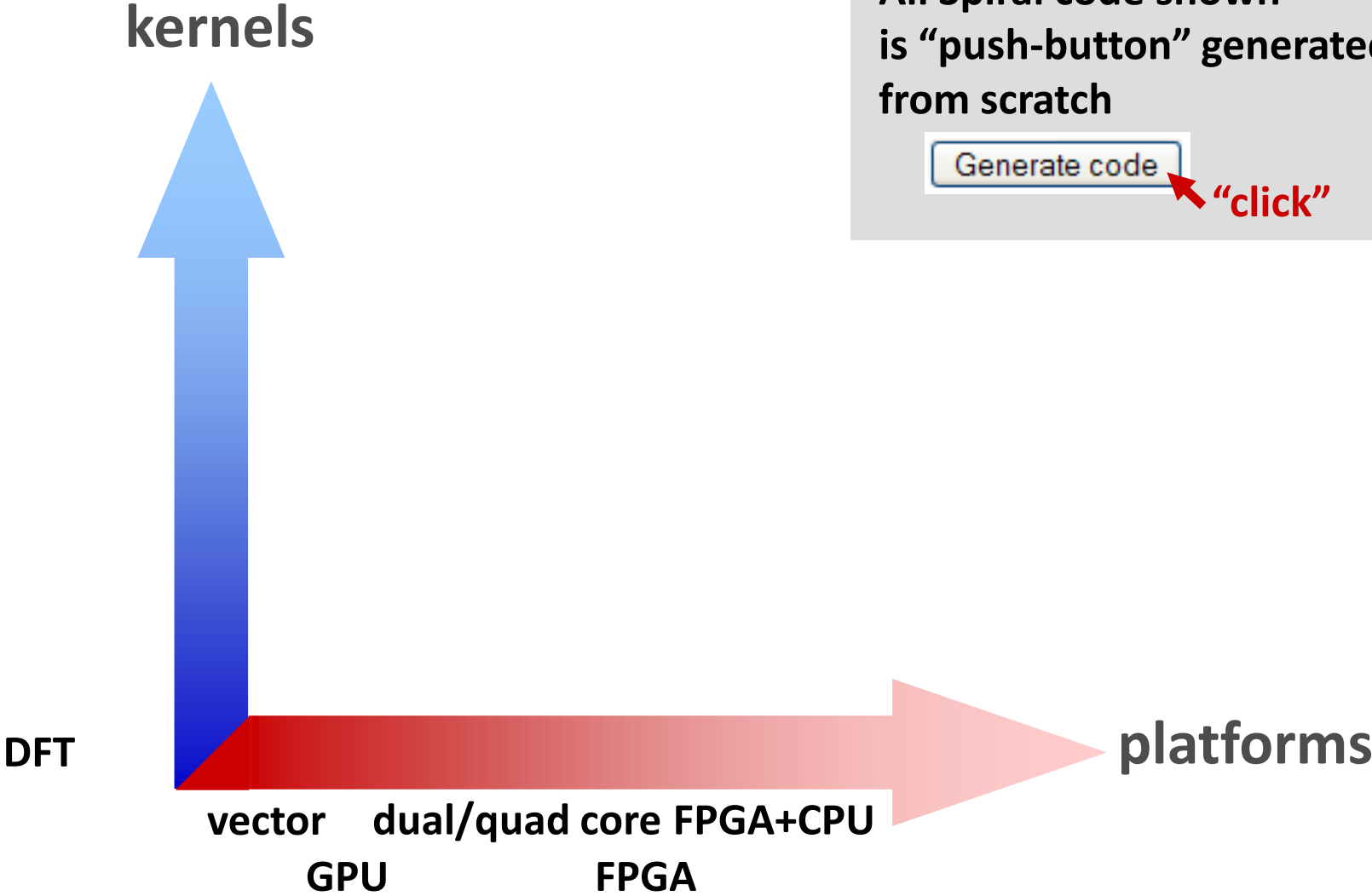
**17  mutually recursive functions**

# Organization

- **Spiral overview**

- **Parallelization in Spiral**

- **Beyond Transforms**

- **Generating general-size libraries**
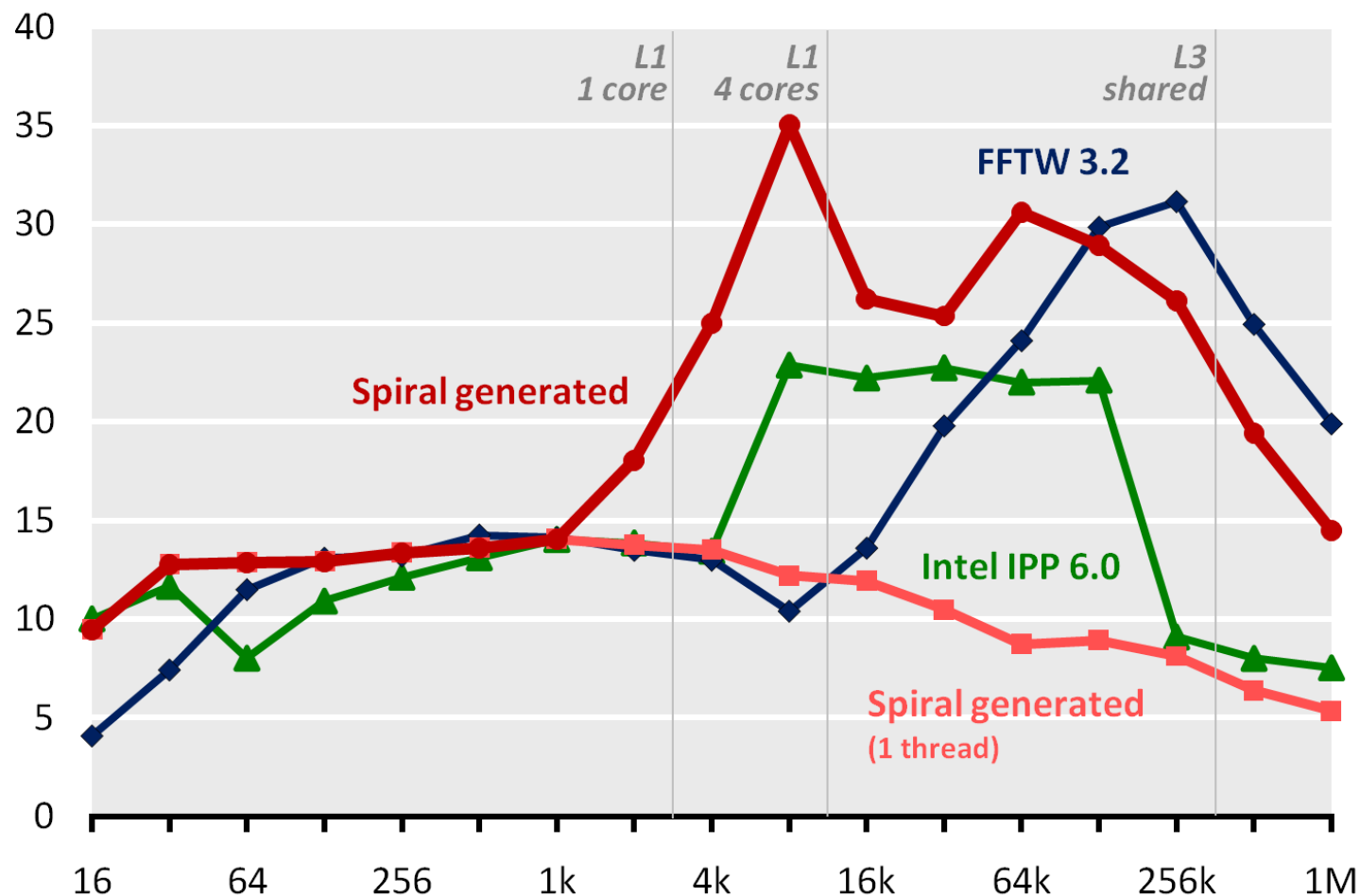
- **Results**

- **Concluding remarks**

# Benchmarks

**kernels**

All Spiral code shown
is "push-button" generated
from scratch

Generate code

"click"

**platforms**

DFT

vector    dual/quad core FPGA+CPU

GPU                 FPGA

# Complex DFT (Intel Core i7, 2.66 GHz, 4 cores, single precision)

performance [Gflop/s] vs. input size

F. Franchetti, M. Püschel: **Short Vector Code Generation for the Discrete Fourier Transform.**
In Proceedings of the 17th International Parallel and Distributed Processing Symposium (IPDPS '03).

F. Franchetti, Y. Voronenko, and M. Püschel: **FFT Program Generation for Shared Memory: SMP and Multicore.**
In Proceedings of Supercomputing, 2006.

# Intel Multicore: Off The Beaten Path



**DCT on 2.66 GHz Core2 (single-precision, 4-way SSSE3)**
performance [Gflop/s], opcount = 2.5 n ld n, Windows XP 32-bit, Intel C++ 10.1

Spiral generated

5–6x

Intel, FFTW

Spiral DCT4
Spiral DCT3
Spiral DCT2
FFTW 3.1.2 DCT4 (e11)
FFTW 3.1.2 DCT3 (e01)
FFTW 3.1.2 DCT2 (e10)
Intel IPP 5.1 DCT

input size n

**DFT on 2.66 GHz Core2 (single-precision, SSSE3)**
performance [Gflop/s], opcount = 5 n ld n, Windows XP 32-bit, Intel C++ 10.1

Spiral x87
Spiral SSE
Intel MKL 9.0

Spiral, Intel

input size

- **DCT: Native algorithm (Spiral) vs. FFT translation (FFTW, MKL)**
  Algorithms developed with the Algebraic Signal Processing theory

- **DFT: SIMD-specific aggressive data layout optimization**
  Included in IPP 6.0 (new domain: IPPGen)

F. Franchetti, M. Püschel:
**SIMD Vectorization of Non-Two-Power Sized FFTs.**
Proceedings of International Conference on Acoustics, Speech, and Signal Processing (ICASSP) 2007.

# Single Node: BlueGene Supercomputers

**DFT, double precision, XL C compiler**

performance [Mflop/s]



BlueGene/L: custom FPU

Legend:
- SPIRAL C99 + 440d
- SPIRAL C + 440d
- SPIRAL C + 440
- FFTW 2.1.5
- GNU GSL

2x

problem size

**DFT, double precision, XL C compiler**

performance [Mflop/s]



BlueGene/P: custom FPU + 4 cores

Legend:
- 4 threads (450d)
- single core (450d)
- single core (450)
- GSL 1.5

3.5x

problem size

**Single BlueGene/L CPU at 700 MHz**
IBM T. J. Watson Research Center

**SIMD vectorization**

**Single BlueGene/P node (4 CPUs) at 850 MHz**
Argonne National Laboratory

**SIMD vectorization + multi-threading**

F. Gygi, E. W. Draeger, M. Schulz, B. R. de Supinski, J. A. Gunnels, V. Austel, J. C. Sexton, F. Franchetti, S. Kral,
C. W. Ueberhuber, J. Lorenz: **Large-Scale Electronic Structure Calculations of High-Z Metals on the BlueGene/L Platform.**
In Proceedings of Supercomputing, 2006. **Winner of the 2006 Gordon Bell Prize (Peak Performance Award).**

J. Lorenz, S. Kral, F. Franchetti, C. W. Ueberhuber: **Vectorization Techniques for the Blue Gene/L double FPU.**
IBM Journal of Research and Development, Vol. 49, No. 2/3, 2005.

SPIRAL
www.spiral.net

# New Multicore Architectures: Cell and GPU



- **Cell BE:** SIMD vector + local stores + DMA

- **GPU:** SIMD vector + OpenGL + Cg + SMT

F. Franchetti, Y. Voronenko, P. A. Milder, S. Chellappa, M. Telgarsky, H. Shen, P. D'Alberto, F. de Mesmay, J. C. Hoe, J. M. F. Moura, M. Püschel: **Domain-Specific Library Generation for Parallel Software and Hardware Platforms.** In Proceedings of The NSF Next Generation Software (NGS) Workshop 2008.
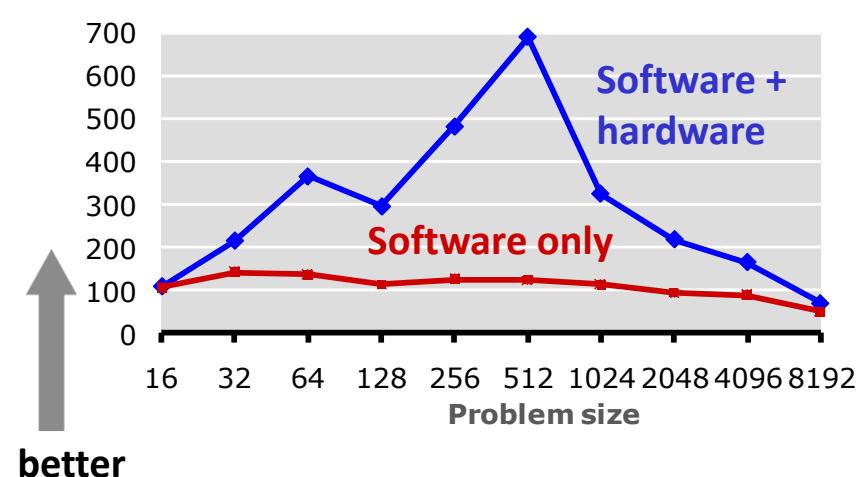
S. Chellappa, F. Franchetti, and M. Püschel: **Computer Generation of Fast FFTs for the Cell Broadband Engine** In Proceedings of the International Conference on Supercomputing (ICS), 2009.

# Hardware: FPGA, CPU + FPGA-Acceleration

**DFT 256 (Verilog Design)**
inverse throughput (gap) [us]



**Spiral**

(Pareto-optimal HW designs)

**Xilinx Logicore 3.2**

Area [slices]

**better** ←

**DFT (CPU accelerated by FPGA)**
performance [Mflop/s]



**Software + hardware**

**Software only**

Problem size

**better**

## Xilinx Virtex 2 Pro FPGA: 1M gates @ 100 MHz + 2 PowerPC 405 @ 300 MHz

P. A. Milder, F. Franchetti, J. C. Hoe, and M. Püschel:
**Formal Datapath Representation and Manipulation for Implementing DSP Transforms.**
In Proceedings of Design Automation Conference (DAC), 2008.

P. D'Alberto, F. Franchetti, P. A. Milder, A. Sandryhaila, J. C. Hoe, J. M. F. Moura, and M. Püschel:
**Generating FPGA Accelerated DFT Libraries.**
In Proceedings of Field-Programmable Custom Computing Machines (FCCM), 2007.

# Qualitative Customization: Code Size



**Performance [Gflop/s]**

- 13 KLOC
- 3 KLOC
- 2 KLOC
- 1.3 KLOC
- 1 KLOC
- FFTW: 150 KLOC

size

Y. Voronenko, F. de Mesmay, and M. Püschel: **Computer generation of general size linear transform libraries**
in Proceedings Code Generation and Optimization (CGO), 2009.

SPIRAL
www.spiral.net

# Benchmarks

**kernels**

All Spiral code shown
is "push-button" generated
from scratch

Generate code    ↖ **"click"**

**GEMM**

**SAR**

**JPEG2000**

**DFT**

**platforms**

vector    dual/quad core FPGA+CPU
GPU              FPGA

# Result: Matrix Multiplication Library

# JPEG 2000: Wavelet + Entropy Coding

**JPEG2000 Image Compression on 2.66 GHz Core2 Duo (2 threads)**
runtime [ms]



**Breakdown:**
**Other**
**Tier 2**
**Tier 1 (EBCOT+MQ)**
**Wavelet (DWT)**

better

**Tile size = image size, resolution level = 2, EBCOT size = 64x64**

**Joint work with S. Hao**

# Polar Format SAR on Intel Core2 Quad

**SAR Image Formation on Intel platforms**

performance [Gflop/s]



Legend:
- 3.0 GHz Core 2 (65nm)
- 3.0 GHz Core 2 (45nm)
- 2.66 GHz Core i7
- 3.0 GHz Core i7 (Virtual)

*newer platforms*

- **Algorithm by J. Rudin (best paper award, HPEC 2007): 30 Gflop/s on Cell**
- **Each implementation: vectorized, threaded, cache tuned, ~13 MB of code**

D. McFarlin, F. Franchetti, M. Püschel, and J. M. F. Moura: **High Performance Synthetic Aperture Radar Image Formation On Commodity Multicore Architectures.** in Proceedings SPIE, 2009.

# Organization

- **Spiral overview**

- **Parallelization in Spiral**

- **Beyond Transforms**

- **Results**

- **Concluding remarks**

# Summary

- **Platforms are powerful yet complicated**
  optimization will stay a hard problem
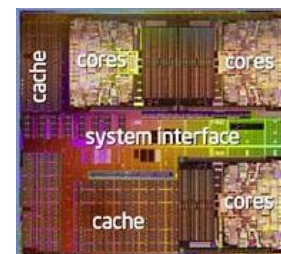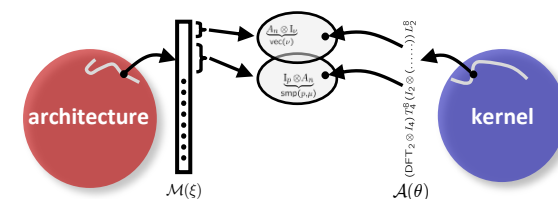


Image: Intel        Image: Intel

- **Unified mathematical framework**
  captures platforms and algorithms

$$\underbrace{I_p \otimes A_n}_{\mathsf{smp}(p,\mu)}$$

- **Program generation and optimization**
  can provide full automation



- **Ongoing Research: Add more functionality**
  *kernels, platforms,...*